

# A Game-Theoretic Approach To Peer Disagreement\*

Remco Heesen<sup>†</sup>      Pieter van der Kolk<sup>‡</sup>

October 21, 2014

## Abstract

In this paper we propose and analyze a game-theoretic model of the epistemology of peer disagreement. In this model, the peers' rationality is evaluated in terms of their probability of ending the disagreement with a true belief. We find that different strategies—in particular, one based on the Steadfast View and one based on the Conciliatory View—are rational depending on the truth-sensitivity of the individuals involved in the disagreement. Interestingly, the Steadfast and the Conciliatory Views can even be rational simultaneously in some circumstances. We tentatively provide some reasons to favor the Conciliatory View in such cases. We argue that the game-theoretic perspective is a fruitful one in this debate, and this fruitfulness has not been exhausted by the present paper.

---

\*Thanks to Kevin Zollman, Jan-Willem Romeijn, Frank Hindriks, Liam Bright, Lauren Leydon-Hardy, and Matt Frise for valuable comments.

<sup>†</sup>Department of Philosophy, Baker Hall 135, Carnegie Mellon University, Pittsburgh, PA 15213-3890, USA. Email: [rheesen@cmu.edu](mailto:rheesen@cmu.edu)

<sup>‡</sup>Faculty of Philosophy, University of Groningen, Oude Boteringestraat 52, 9712 GL Groningen, The Netherlands. Email: [p.m.van.der.kolk@rug.nl](mailto:p.m.van.der.kolk@rug.nl)

# 1 Introduction

The aim of this paper is to show that the problem of peer disagreement can be analyzed from a game-theoretic perspective. The problem of peer disagreement, as it is presented in the literature, is how to respond rationally to the disagreement from an epistemic peer,<sup>1</sup> whereby *epistemic peer* is construed as an agent who has the same evidence and is comparably good at evaluating that evidence.<sup>2</sup> Game theory, in turn, is the study of strategic decision making, where “strategic” means that the decision of one decision maker may interact with that of another. This paper explains how the latter can be used to analyze the former.

To do so, we focus on two prominent recommended strategies in the literature about peer disagreement, namely the response advocated by the *Conciliatory View* and the one suggested by the *Steadfast View*.<sup>3</sup> On the Conciliatory View, it can *not* be rational for an agent to stick to her opinion when it is disputed by an epistemic peer. Instead, she should suspend judgment (Feldman 2007), split the difference (Elga 2007), or at least migrate her opinion significantly in the direction of her peer’s conflicting opinion. (Christensen 2007). According to the Steadfast View, on the other hand, it *can* be rational for an agent to retain her opinion in the face of peer disagreement (Kelly 2005, van Inwagen 2010).

The game-theoretic toolkit enables us to analyze the rationality of these responses (strategies) for disagreeing peers (players), relative to these peers’ epistemic goals (preferences). In the literature on peer disagreement, the epistemic goal or preference is commonly understood to be believing the

---

<sup>1</sup>See, for example, Kelly (2005, 167), Christensen (2009, 756), Elga (2007, 478), and Feldman (2007, 201). For more papers about the issue, see the collection of papers in Feldman and Warfield (2010), and Lackey and Christensen (2013), and papers cited below.

<sup>2</sup>For such construals of peerhood, see, for example, Kelly (2005, 170), Christensen (2007, 188), Feldman (2007, 201), and Lackey (2008, 274).

<sup>3</sup>In the debate about peer disagreement, it is common to talk about “responses”, whereas in the context of game theory “strategies” is conventional. In this paper, we will use the two terms interchangeably.

correct truth-value of the proposition under discussion.<sup>4</sup> Thus, the rationality of the available responses—i.e., the Conciliatory strategy and the Steadfast strategy—can be analyzed by investigating to what extent they satisfy the preferences (epistemic goals) of the disagreeing peers. In section 2 and 3 we will further explain the details of this game-theoretic approach to the problem of peer disagreement. Section 4 discusses the results of this model, section 5 considers some possible extensions or variations of the model, and section 6 wraps up by emphasizing some key take-aways.

Why should such a game-theoretic analysis be a relevant contribution to the debate about peer disagreement? Our motivation is that the resources of game theory enable a clarification of the rationality of peer disagreement—in particular, of the Conciliatory View and the Steadfast View—along an independently motivated and well-developed standard. In the debate about peer disagreement, it is not always clear how exactly rationality is understood, what exactly counts as a peer, what a disagreement is, or even what the Conciliatory View and the Steadfast View exactly amount to.<sup>5</sup> A formalization along the lines of game theory forces us to be precise about these notions and disclose their exact specifications. And the fruit of such explicitness is that it helps us to gain a better understanding of the exact specifications and conditions under which a particular strategy (like the ones suggested by the Conciliatory View and the Steadfast View) can be considered a rational response to the disagreement from a peer.

We do not want to suggest that our game-theoretic model is the best, let alone the only, way to make the machinery under the problem of peer disagreement formally precise. Rather, our aim is to show that it *can* be done. And we would welcome different or variant precisifications, as we think that this would only help the debate.

---

<sup>4</sup>See, for example, Christensen (2007, 216), Feldman (2007, 212), Elga (2007, 488), Kelly (2010, 17), and, even if only indirectly, White (2005, 450).

<sup>5</sup>This complaint is also expressed by, for example, Fitelson and Jehle (2009) Moss (2011), and Lasonen-Aarnio (2013).

## 2 The Peer Disagreement Game

We introduce our game-theoretic setup with the help of an informal example. Imagine two detectives, call them Jane and Hercule, who both have been asked to go to a crime scene to investigate whether  $\phi$ , say, whether the butler is the culprit. We make the following three assumptions about the detectives. First, they have the same evidence at their disposal to investigate  $\phi$ , namely whatever traces are left at the crime scene. Second, the detectives can make an informed estimation of how reliable each of them is in investigating  $\phi$ , based on their respective *track-records*; the number of crimes they have solved in the past compared to the number of crimes they didn't solve. Third, the detectives really want to find out the truth regarding  $\phi$ , they really want to solve the case.<sup>6</sup>

So, Jane and Hercule both go to the crime scene, and spend some time examining and evaluating the evidence. After some time, they meet up to report their findings.

In the context of our “peer disagreement game”, two things can happen at this point. Jane and Hercule can either have formed the same belief about  $\phi$ , or they can have formed conflicting beliefs and disagree about  $\phi$ .<sup>7</sup>

If the detectives have reached the same conclusion about  $\phi$ , say, they agree that the butler is indeed the culprit, then there is no problem of peer disagreement. The detectives can go write their reports. The case that we are interested in, of course, is when the detectives have formed conflicting opinions regarding  $\phi$ ; for example, when Jane believes that the detective is the culprit and Hercule believes that the detective is innocent. And our question is what, in such a case, a rational response for Jane and Hercule can

---

<sup>6</sup>We take it that the fulfilment of these three conditions is what is (at minimum) required for the two detectives to be called each other's peers, considering the construals of peerhood by, for example, Kelly (2005, 175), Elga (2007, 484), Lackey (2008, 274), and Christensen (2009, 757). The attribution of peerhood then depends on how equal the detectives must be in their reliability. Our analysis accommodates this.

<sup>7</sup>In this paper we restrict ourselves to full belief states. An analysis that includes degrees of belief is possible (see section 5), but we leave it to future work.

be, given their goal of finding out the truth about  $\phi$ , and the information they have about each other's track-records.

Based on the debate about peer disagreement, we distinguish three strategies that the detectives can play. The first comes from the Steadfast View and is the strategy of staying with the initial belief. We call this strategy **Stay**. The second strategy is the Conciliatory View's recommendation to suspend judgment.<sup>8</sup> This strategy is called **Suspend**. And third, for the sake of completeness, we include switching to the belief of the other detective as a third possible strategy, called **Switch**.

So, after Jane and Hercule find out that they disagree about whether the butler is the culprit, they can each play one of these three strategies. When, say, Jane plays **Suspend**, she withdraws her initial belief about  $\phi$ , goes back to the crime scene to re-examine the evidence, and forms a new belief about  $\phi$ . But when Jane plays **Stay**, she chooses to ignore the disagreement from Hercule and maintains her initial opinion, no matter what Hercule does. And when Jane plays **Switch**, she chooses to ignore her own opinion and takes over the belief of Hercule, regardless of what that belief is.<sup>9</sup>

The disagreement game ends when the two detectives reach an agreement about  $\phi$ . For example, when Jane believes that the butler did it, and Hercule

---

<sup>8</sup>Technically, the recommendation from the Conciliatory View can also be to split the difference with the opponent (Elga 2007), or to revise one's initial confidence level in the proposition considerably (Christensen 2007). But since we restrict ourselves to full beliefs, we take it that the Conciliatory View's recommendation amounts to suspending judgment.

<sup>9</sup>So, naturally, only when a detective plays **Suspend** she gets a chance at forming a new opinion. It might be objected that acquiring a new belief is not a necessary consequence of suspending judgment. We agree. We think that it is in the spirit of suspending judgment, in cases of peer disagreement (or other significant counterevidence), that judgment is suspended only *momentarily*, as an act of caution, to re-examine the evidence and check whether the initial opinion was correct. But, on another reading, one might say that suspending judgment is meant for the long run, or at least until new evidence comes in, precisely because forming any belief about the matter would be irrational. In this paper we work with the first reading. But the second reading might be another welcome extension of our analysis (see section 5).

believes that he didn't, and Jane plays **Suspend** and Hercule plays **Stay**. Then the game ends when, after re-examining the evidence, Jane draws the same conclusion as the conclusion that Hercule was holding on to, namely that the butler is innocent. The same would happen when, for example, Jane would play **Stay** and Hercule would play **Switch**. But the game continues when, after playing their strategy, the two detectives still disagree about  $\phi$ . For example, when the detective playing **Suspend** forms a new belief about  $\phi$  that again conflicts with the conclusion of the other detective, then this detective will again suspend judgment, re-examine the evidence, and come back with a new belief.

For the purposes of this paper, we assume that the detectives do not change strategies throughout the disagreement game.<sup>10</sup> This means that the game might also continue forever. For example, when Jane believes that the butler is innocent and Hercule disagrees, and both detectives play **Stay**, then they will never come to an agreement. The same thing happens when both detectives play **Switch**.<sup>11</sup>

And now we are in a position to analyze how well these strategies do—in particular, the Conciliatory strategy **Suspend** and the Steadfast strategy **Stay**—in guiding each detective to the correct verdict on whether the butler did it. Which of these strategies gives a detective the best prospects of arriving at the truth?

Observe that which strategy is best is going to depend on two factors. First, it depends on the track-record, the reliability of each of the two detectives. For example, if Jane thinks that Hercule is way better at evaluating

---

<sup>10</sup>The reason is that this allows a straightforward comparison of the Conciliatory View, which recommends playing **Suspend** for all instances of peer disagreement, and the Steadfast View, according to which playing **Stay** can be rational. It would be an interesting extension of our model to allow players to change their strategy during the game (see section 5).

<sup>11</sup>That under these strategies the game continues forever doesn't make an evaluation of the rationality of these strategies impossible. For in both cases we can still evaluate how well these strategies do with respect to tracking the truth.

correctly whether the butler did it, then it would be ill-advised for her to play **Stay** upon finding out that Hercule disagrees with her initial assessment. But when Jane thinks that she is far more reliable than Hercule, then playing **Stay** is quite sensible.

Second, which strategy is best depends also on the strategy of the other player. For example, when Hercule plays **Stay**, then it does not really matter for Jane whether she plays **Suspend** or **Switch**, because either way the game will end when Jane takes over the conclusion of Hercule. But when Hercule plays **Switch**, it *does* matter whether Jane plays **Suspend** or **Switch**, because playing **Switch** will bring them in a state of perpetual disagreement, whereas playing **Suspend** will make them agree eventually. We will return to these points in section 4.

So, we are going to analyze when a strategy played by a detective—say, the Conciliatory **Suspend**, or the Steadfast **Stay**—is the best possible strategy to arrive at the truth about the butler, when taking into account the strategy played by the other detective, as well as the detective’s own reliability and the reliability of the other detective.

This concludes our informal description of the peer disagreement game. In the next section we will provide the formal vocabulary, and then analyze this game.

### 3 Rationality for Jane and Hercule

Whenever Jane and Hercule investigate the evidence, they may conclude that the butler did it ( $\phi$ ) or that he did not do it ( $\neg\phi$ ). One of these conclusions is *true* and one is *false*.

We will denote by  $p$  and  $q$  the reliability or *truth-sensitivity* of Jane and Hercule, respectively. Thus  $p$  is the probability, on any given investigation, that Jane draws a true conclusion from the evidence.  $1 - p$  denotes the probability of a false belief. So if the butler really did it Jane believes that he did it with probability  $p$  and believes in his innocence with probability  $1 - p$ .

Whereas if he is innocent she believes in his innocence with probability  $p$  and believes that he did it with probability  $1 - p$ . Hercule’s probabilities of drawing a true or a false conclusion from the evidence are denoted by  $q$  and  $1 - q$ , respectively.<sup>12</sup>

We choose to model “the probability of generating a true or false belief” rather than “the probability of generating a belief for or against  $\phi$ ” because we have evidence for the former but not the latter based on the respective track-records of the two detectives. We assumed at the start of section 2 that this track-record information is known to the two detectives.

In the epistemology of peer disagreement—as we learn from, for example, Christensen (2007, 216), Feldman (2007, 212), Elga (2007, 488), and Kelly (2010, 17)—the objective of rational conduct is commonly understood to be believing the correct truth-value. This suggests the following epistemic norm.

**Truth Norm (TN).** Having a true belief is more valuable than having a false belief.

We assume that Jane and Hercule share this noble goal, and that in fact obtaining a true belief about whether the butler did it is their *only* goal.<sup>13</sup> So the two detectives are not distracted by pragmatic concerns. This is a methodological, not a substantive assumption: we are interested in the epistemology of peer disagreement, not its pragmatics.

Then we can easily model their *preference* over *outcomes* of the disagreement game: Jane prefers an outcome in which she has a true belief about the butler’s guilt over one in which she has a false belief, and likewise for

---

<sup>12</sup>To avoid trivial cases, we assume that  $0 < p < 1$  and  $0 < q < 1$ . We assume that, if Jane or Hercule suspends judgment in response to disagreement, their new opinion is generated with the same probabilities as their initial opinion (so Jane believes correctly with probability  $p$ , and Hercule believes correctly with probability  $q$ ). We also assume that each time an opinion is generated this is done independently (in the probabilistic sense) from the detective’s previous opinions and the other detective’s current or previous opinions.

<sup>13</sup>We recognize that one might have other epistemic goals than truth. We will discuss this in section 5.



Hercule.<sup>14</sup> A detective receives utility 1 if his/her belief about the guilt or innocence of the butler at the end of the disagreement game is true, and utility 0 if it is false.<sup>15</sup>

The expected utility of a detective in the game is then simply the probability of ending the game with a true belief. So Jane and Hercule prefer a strategy if it increases their probability of ending the disagreement game with a true belief concerning  $\phi$ .

We can now determine the probabilities of ending the disagreement game with a true belief for each combination of strategies of the two players (a combination of strategies is called a *strategy profile*).

If both detectives play **Stay**, they never change their mind in response to disagreement, so their probability of ending with a true belief is simply the probability that they obtain a true belief initially:  $p$  for Jane and  $q$  for Hercule. In all other cases the probability of ending the disagreement game with a true belief is the same for both detectives. These probabilities are indicated in table 1. The rows of table 1 indicate Jane's choice of strategy, and the columns indicate Hercule's choice.<sup>16</sup>

---

<sup>14</sup>Note that under our interpretation of (TN) players care only about the truth of their own belief. Results concerning a variation of our model where players also care about the truth of the other player's belief are available from the authors upon request.

<sup>15</sup>The introduction of utilities here adds nothing over and above the informal statement in the previous sentence. In particular the numbers 0 and 1 are arbitrary: all that matters is that a true belief yields a higher utility.

<sup>16</sup>This completes our specification of the game. Formally, a game is a triple  $(N, \{S_i\}_{i \in N}, \{u_i\}_{i \in N})$ , where  $N$  is the set of players,  $S_i$  the set of strategies available to player  $i$ , and  $u_i$  the utility function for player  $i$ , which assigns real-valued utility to each strategy profile. In our case there are two players:  $N = \{\text{Jane}, \text{Hercule}\}$ ; the strategy sets for both players are identical:  $S_{\text{Jane}} = S_{\text{Hercule}} = \{\text{Stay}, \text{Suspend}, \text{Switch}\}$ ; and the utility for each player on each strategy profile is as in table 1.

The utilities are determined using the description of the disagreement game given in section 2. For example, if both detectives play **Suspend** they will generate new beliefs repeatedly until the first time they agree. The probability that they both generate a belief that  $\phi$  is true is  $pq$  and the probability that they agree that  $\phi$  is false is  $(1 - p)(1 - q)$ . So the probability that they end the game with a correct belief about  $\phi$  is the probability that, on the first round on which they agree, they agree that  $\phi$  is true rather than that  $\phi$

	Stay	Suspend	Switch
Stay	$(p, q)$	$p$	$p$
Suspend	$q$	$\frac{pq}{pq+(1-p)(1-q)}$	$\frac{p(1-(1-p)(1-q))}{1-p(1-p)}$
Switch	$q$	$\frac{q(1-(1-p)(1-q))}{1-q(1-q)}$	$pq$

Table 1: Expected utilities associated with each strategy profile under (TN).

How can the detectives maximize their probability of ending the disagreement game with a true belief, given that the choice of strategy of the other detective influences his or her probability of attaining true belief, but they cannot control it? Game theorists have invented various concepts of rationality in a game to deal with this problem. We will use the notion of Nash equilibrium.

A *Nash equilibrium* is a profile—that is, an assignment of a strategy to each player—in which either player’s strategy is a best response to the other’s. In other words, in a Nash equilibrium, no player can get an outcome she prefers over the equilibrium outcome by unilaterally changing her strategy. In our game this means that in a Nash equilibrium Jane and Hercule are maximizing their respective probabilities of ending the game with a true belief, *given* (that is, keeping fixed) the other detective’s strategy. This is how we interpret (epistemic) rationality for Jane and Hercule.<sup>17</sup>

---

is false. This probability is simply  $pq$  divided by  $pq + (1 - p)(1 - q)$ .

<sup>17</sup>Our model is thus different from other game-theoretic models that deal with information and disagreement, notably Aumann (1976). There are at least three important differences. First, there is an iterated exchange of opinions rather than common knowledge of posteriors, although Geanakoplos and Polemarchakis (1982) show how Aumann’s result may be reinterpreted in an iterative context. Second, we work with full belief states rather than degrees of belief. Third, we allow players to choose a strategy for updating their beliefs, rather than assuming Bayesian updating. We do not wish to argue that our assumptions necessarily yield a better model of peer disagreement, in fact we would welcome a comparison of our model with Aumann’s in terms of their consequences for the peer disagreement debate.

## 4 Results and Discussion

What are the Nash equilibria of this game?<sup>18</sup> This turns out to depend on the values of  $p$  and  $q$ . Figure 1 shows which strategy profiles are Nash equilibria for any combination of values of  $p$  and  $q$ .

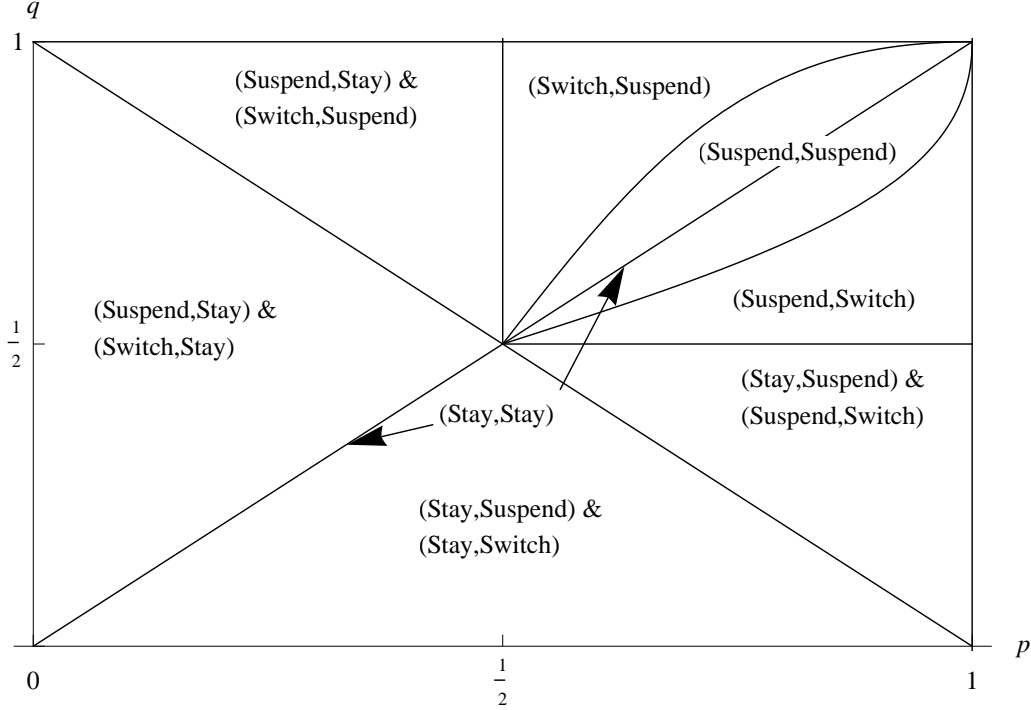


Figure 1: Nash equilibria of the peer disagreement game as a function of the truth-sensitivity of the detectives. E.g., if  $p = 0.4$  and  $q = 0.9$  then the Nash equilibria of the game are **(Suspend, Stay)** and **(Switch, Suspend)**.

Recall that we noted in section 2 that two factors would influence which strategy choice is best. First, the reliability of the two detectives (modeled as  $p$  and  $q$ ) and second, the strategy of the other detective. Both of these factors are shown clearly in our results in figure 1.

<sup>18</sup>We consider only pure strategy equilibria.

The values of  $p$  and  $q$  clearly influence which strategy profiles are rational. For example, **(Stay,Switch)** is a Nash equilibrium whenever  $q \leq \min\{p, 1 - p\}$ , but not otherwise. Similarly, **(Switch,Stay)** is a Nash equilibrium whenever  $p \leq q \leq 1 - p$ , but not otherwise.

The other detective's strategy also clearly influences what it is rational for a detective to do. For example, if  $1 - p \leq q \leq 1/2$  the Nash equilibria are **(Stay,Suspend)** and **(Suspend,Switch)**. So under these circumstances, if Hercule chooses the strategy **Suspend**, it is rational for Jane to choose **Stay**, but if Hercule chooses **Switch**, it is rational for Jane to choose **Suspend**.

The epistemic success of the two detectives (both in terms of which strategy promises the best probability of a true belief, and in terms of the value of that probability) thus depends heavily on the choices made by the other detective. In this way the epistemology of this model is truly *social*.

Of particular interest in evaluating the results in figure 1 are the profiles **(Stay,Stay)** and **(Suspend,Suspend)**. This is because the former captures most directly the Steadfast View — according to which it can be rational to **Stay** in a case of peer disagreement — and the latter captures most directly the Conciliatory View — according to which the only rational option is to **Suspend**.

Perhaps surprisingly, *both* **(Stay,Stay)** and **(Suspend,Suspend)** turn out to constitute Nash equilibria, under some conditions even both at once.

As we can see from figure 1, the Steadfast profile **(Stay,Stay)** is a Nash equilibrium when Jane and Hercule are each other's equals in terms of how truth-sensitive their beliefs are (i.e.,  $p = q$ ). In such a case neither would gain anything by playing **Suspend** or **Switch** (provided the other player continues to play **Stay**). More precisely, when  $p = q$ , the probability that a detective ends up with a true belief by staying with his or her initial opinion is just as high as the probability that the opinion of the other detective or a newly generated opinion is true.

However, a mutual Conciliatory approach, as expressed in the strategy profile **(Suspend,Suspend)**, can *also* be a Nash equilibrium. This happens

when  $p$  and  $q$  are both greater than  $1/2$  and are relatively close to each other (see figure 1).<sup>19</sup> When both detectives are relatively good learners, and they find out that they have formed conflicting beliefs, they stand to gain more when they both suspend judgment and acquire a new belief, than when they stick to their initial beliefs, or switch to the other detective's belief.

An especially interesting scenario occurs whenever  $p$  and  $q$  are exactly equal and greater than  $1/2$ : then (Stay,Stay) and (Suspend,Suspend) are Nash equilibria at the same time. Under the definition of rationality we use, in such a case both Steadfast and Conciliatory strategies are rational. However, we offer three reasons to think that the Conciliatory strategy should be preferred (without necessarily endorsing these reasons as decisive).

First, whenever (Stay,Stay) and (Suspend,Suspend) are Nash equilibria simultaneously, (Suspend,Suspend) offers a higher utility (a higher probability of solving the case correctly) to both detectives.<sup>20</sup> In fact, (Suspend,Suspend) is *Pareto efficient*. So Jane and Hercule prefer to play (Suspend,Suspend) over (Stay,Stay). If they are allowed to discuss their strategy before the game starts, we should expect both detectives to play **Suspend**.

Second, **Suspend** is a *weakly dominant* strategy (for both detectives), while **Stay** is not. This means that playing **Suspend** pays off at least as well as playing **Stay** or **Switch**, *regardless* of what strategy the other detective chooses. So in this situation, playing **Stay** is only best for a detective who is absolutely certain that the other detective is playing **Stay** as well (and even then playing **Suspend** is equally good), whereas if there is only the slightest uncertainty about what the other detective is going to do, **Suspend** is the uniquely best strategy.

Third, we can see in figure 1 that when  $p$  and  $q$  are both greater than  $1/2$  there is a significant area in which the profile (Suspend,Suspend) is a Nash

---

<sup>19</sup>More precisely, the region where (Suspend,Suspend) is a Nash equilibrium is characterized by the inequality  $\frac{p-\sqrt{p(1-p)}}{2p-1} \leq q \leq \frac{p^2}{1-2p(1-p)}$  (although the first expression is undefined when  $p = 1/2$ , the point  $p = q = 1/2$  is also part of this region).

<sup>20</sup>Whenever  $p = q > 1/2$ , it must also be the case that  $\frac{p^2}{p^2+(1-p)^2} > p$ .

equilibrium, while (Stay,Stay) is a Nash equilibrium only when  $p$  and  $q$  are exactly equal. This means that the strategy **Suspend** has a larger margin for error than the strategy **Stay**. If Jane and Hercule lack precise information about each other’s truth-sensitivity (as is reasonable to expect), playing **Stay** is “riskier” than playing **Suspend** because the former requires exact and the latter only approximate equality of the player’s truth-sensitivities.

To sum up, a surprising result of this model is that if the detectives are both good learners, and the one is not significantly better than the other, then both the Steadfast strategy and the Conciliatory strategy can be Nash equilibria. However, we have noted three reasons to think that in such cases the Conciliatory strategy should be preferred.

## 5 Limitations and Extensions of our Analysis

We have limited our analysis to a particular game-theoretic formalization of a particular disagreement game between two detectives, Jane and Hercule. To what extent does our analysis generalize to other peer disagreements? And what are possible variations or extensions of our formalization?

Regarding the first question, our analysis applies to peer disagreements in general insofar as they satisfy the assumptions of our model. In particular, (1) peers are cashed out in terms of comparable reliability or truth-sensitivity, (2) the possible responses available to the peers are something like the strategies **Stay**, **Suspend**, and **Switch** as we model them, and (3) the rationality of a particular response is evaluated in terms of how well it tracks the truth.

Regarding the second question, there are many options for different peer disagreement games. Let us give six variables that can be filled in differently.

First, *doxastic attitudes*: in our model, strategies act on full belief states, mainly for reasons of simplicity, but strategies might also be interpreted as adjusting degrees of belief.

Second, we forced our detectives to generate a new belief whenever they suspend judgment on  $\phi$ . A variation of our model might allow peers to

persist in a state of suspension. This outcome could be assigned its own value, presumably worse than having a true belief but better than having a false belief.

Third, we kept the peers' strategies fixed throughout the game. The reason for this was to enable an evaluation of the Conciliatory and Steadfast strategies. But it is of course an idealization. So it would be an interesting extension of the game to allow peers to change their strategies during the game.

Similarly, we assumed that the game might go on indefinitely. This is not very realistic. In real life there are time and energy constraints. So another possible extension would be to let the game continue for a limited number of rounds, after which the agents must have made up their minds.

Fifth, in our analysis the rationality of a strategy was evaluated using Nash equilibria. Although this is very natural in game theory, it has substantive normative implications. So one may want to consider alternatives. Available alternatives include various refinements of the notion of equilibrium, such as the trembling hand equilibrium, and alternative standards, such as weak dominance. Different strategies may turn out to be rational under such different standards of rationality.

Finally, we worked with only one epistemic norm, namely truth. But there are more epistemic goals. For example, many philosophers of science have argued, under the label of "epistemic diversity", that maintaining diversity of opinion can have epistemic value to a population of scientists, stimulating new ideas and discoveries (Feyerabend 1975, Kitcher 1990, Zollman 2010). Similarly, the literature on epistemic rationality has identified a trade-off between truth and information (Levi 1967). For example, true beliefs could be maximized by believing only tautologies, but this is not informative. Either of these considerations could motivate augmenting or replacing (TN) with different norms.

## 6 Conclusion

By way of conclusion we would like to emphasize four lessons that can be drawn from our preliminary game-theoretic investigation of the epistemology of peer disagreement.

First, in our model the Steadfast and Conciliatory Views were sometimes both right: there were circumstances in which both staying with your own opinion and suspending belief were rational. The idea that staying and suspending can be rational simultaneously is underexplored in the literature and worth investigating more extensively.

Second, the rationality of a response to peer disagreement may depend on the truth-sensitivity of the peers. Both the peers' relative truth-sensitivity (who is a better learner and by how much?) and their absolute truth-sensitivity (are they better than chance, say, or some other objective threshold?) can make a difference.

Third, what is rational for a peer to do (e.g., whether to be Steadfast or Conciliatory) may depend on what the other peer is doing. This is a natural conclusion to draw in the game-theoretic context, but underexplored in the peer disagreement literature.

Fourth, analysis of other game-theoretic models of peer disagreement may shed more light on the above three points and other important questions about peer disagreement. We encourage anyone interested in our model (especially if they liked it but for one or two assumptions) to develop and analyze such an alternative game-theoretic model of peer disagreement. We hope to have provided a fruitful framework within which such further models can be developed.

## References

Robert J. Aumann. Agreeing to disagree. *The Annals of Statistics*, 4(6): 1236–1239, 1976.



- David Christensen. Epistemology of disagreement: The good news. *The Philosophical Review*, 116(2):187–217, 2007.
- David Christensen. Disagreement as evidence: The epistemology of controversy. *Philosophy Compass*, 4(5):756–767, 2009.
- Adam Elga. Reflection and disagreement. *Noûs*, 41(3):478–502, 2007.
- Richard Feldman. Reasonable religious disagreements. In Louise Antony, editor, *Philosophers Without Gods: Meditations on Atheism and the Secular*, pages 194–214. Oxford University Press, 2007.
- Richard Feldman and Ted Warfield, editors. *Disagreement*. Oxford University Press, 2010.
- Paul Feyerabend. *Against Method*. New Left Books, London, 1975.
- Branden Fitelson and David Jehle. What is the “Equal Weight View”? *Episteme*, 6(3):280–293, 2009.
- John D. Geanakoplos and Heraklis M. Polemarchakis. We can’t disagree forever. *Journal of Economic Theory*, 28(1):192–200, 1982.
- Thomas Kelly. The epistemic significance of disagreement. In Tamar Szabó Gendler and John Hawthorne, editors, *Oxford Studies in Epistemology*, volume 1, pages 167–196. Oxford University Press, 2005.
- Thomas Kelly. Peer disagreement and higher-order evidence. In Richard Feldman and Ted Warfield, editors, *Disagreement*, pages 111–174. Oxford University Press, 2010.
- Philip Kitcher. The division of cognitive labor. *The Journal of Philosophy*, 87(1):5–22, 1990.
- Jennifer Lackey. What should we do when we disagree? In Tamar Szabó Gendler and John Hawthorne, editors, *Oxford Studies in Epistemology*, volume 3, pages 274–93. Oxford University Press, 2008.

- Jennifer Lackey and David Christensen, editors. *The Epistemology of Disagreement: New Essays*. Oxford University Press, 2013.
- Maria Lasonen-Aarnio. Disagreement and evidential attenuation. *Noûs*, 47 (4):767–794, 2013.
- Isaac Levi. *Gambling with Truth*. MIT Press, Cambridge, 1967.
- Sarah Moss. Scoring rules and epistemic compromise. *Mind*, 120(480):1053–1069, 2011.
- Peter van Inwagen. We’re right. they’re wrong. In Richard Feldman and Ted Warfield, editors, *Disagreement*, pages 10–28. Oxford University Press, 2010.
- Roger White. Epistemic permissiveness. *Philosophical Perspectives*, 19(1): 445–459, 2005.
- Kevin J. S. Zollman. The epistemic benefit of transient diversity. *Erkenntnis*, 72(1):17–35, 2010.